

PHISHING ATTACK DETECTION USING MACHINE LEARNING METHOD

JOHN ARTHUR JUPIN

Bachelor of Computer Science
(Computer Systems and Networking) with Honors

UNIVERSITI MALAYSIA PAHANG

SUPERVISOR'S DECLARATION

I hereby declare that I have checked this thesis, and, in my opinion, this thesis is adequate in terms of scope and quality for the award of Bachelor of Computer Science (Computer System and Networking).

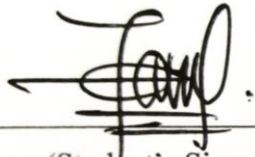


(Supervisor's Signature)

Full Name : DR. MOHD ARFIAN BIN ISMAIL
Position : SENIOR LECTURER
Date : 8 JANUARY 2019

STUDENT'S DECLARATION

I hereby declare that the work in this thesis is based on my original work except for quotations and citations which have been duly acknowledged. I also declare that it has not been previously or concurrently submitted for any other degree at Universiti Malaysia Pahang or any other institutions.



(Student's Signature)

Full Name : JOHN ARTHUR JUPIN

ID Number : CA15069

Date : 8 JANUARY 2019

PHISHING ATTACK DETECTION USING
MACHINE LEARNING METHOD

JOHN ARTHUR JUPIN

Thesis submitted in fulfillment of the requirements
for the award of the degree of
Bachelor of Computer Science (Computer Systems and Networking) with Honors

Faculty of Computer Science & Software Engineering
UNIVERSITI MALAYSIA PAHANG

JANUARY 2019

ACKNOWLEDGEMENTS

First, I would like to thank to my supervisor, Dr. Mohd Arfian bin Ismail, Senior Lecturer of Faculty of Computer Science and Software Engineering of University Malaysia Pahang for helping me from the start of this research until the writing of the thesis for this research. He had given me his full authority on my research yet monitor my research progress. He also had given his advices on how this research should be conducted.

I also would like to thank Prof. Dr. Kamal Zuhairi Zamli as the dean of Faculty of Computer Science and Software Engineering of University Malaysia Pahang for giving students the chance to use the lab during non-academic session.

Finally, I would like to thank my family members, friends and my lectures for supporting and guiding me from the start of this research until the end of this research.

ABSTRAK

Pembangunan rangkaian komputer pada masa kini adalah sangat pesat. Perkara ini tidak dapat lagi dinafikan kerana setiap pengguna komputer di serata dunia perlu menyambungkan komputer mereka ke rangkaian Internet. Hal ini demikian menunjukkan bahawa penggunaan rangkaian Internet adalah sangat penting, sama ada ianya digunakan untuk tujuan kerja dan tugas mahupun untuk mengakses ke akaun media sosial, contohnya Instagram, Facebook dan Twitter. Walaubagaimanapun, dalam penggunaan luas rangkaian komputer ini, secara tidak langsung, privasi pengguna komputer adalah dalam bahaya. Hal ini demikian disebabkan pengguna komputer tidak menitikberatkan sistem sekuriti dan keselamatan di dalam komputer mereka. Oleh sebab ini, penggodam akan menggodam dan membuat serangan rangkaian ke atas pengguna komputer dengan mudah. Hal ini demikian sangat bahaya, terutama sekali kepada organisasi penting kerana penggodam dapat melumpuhkan sistem atas talian dalam syarikat, mencuri maklumat-maklumat sulit dan seterusnya mencuri wang syarikat secara atas talian tanpa disedari oleh mana-mana pihak. Antara serangan yang boleh dibuat termasuklah serangan penafian-perkhidmatan, serangan perdayaan dan *phishing*. Matlamat dalam kajian ini ialah untuk menggunakan alat *anti-phishing* dalam menghalang serangan sekuriti rangkaian di dalam satu organisasi. Dalam kajian ini, serangan *phishing* telah dikaji secara mendalam. Selepas kajian dibuat, cara yang telah diutarakan untuk menghalang serangan *phishing* ialah melalui pembelajaran mesin. Selain itu, kajian ini juga menunjukkan bahawa serangan *phishing* selalunya berhubung kait dengan serangan mesej *spam*. Mesej *spam* ini termasuklah email dan juga mesej SMS yang diterima dari pengguna. Dalam pembelajaran mesin, terdapat beberapa algorithma yang boleh digunakan dalam menghalang kedua-dua serangan ini. Algorithma *Naïve Bayes*, algorithma Pokok Keputusan dan algorithma Mesin Vektor Sokongan telah digunakan untuk menghalang serangan *spam*, dan juga serangan *phishing* daripada berlaku. Kajian algorithma ini dibuat secara mendalam dan cara-cara dalam melaksanakan algorithma ini telah dibincangkan secara mendalam dan lebih terperinci. Eksperiment juga telah dijalankan untuk set data yang diperoleh dengan menggunakan kaedah pembelajaran mesin. Keputusan telah diperoleh, di mana menunjukkan prestasi untuk kaedah pembelajaran mesin untuk setiap set data.

ABSTRACT

The development of computer networks today is increased rapidly. This can be shown based on the trend of every computer user around the world, whereby they need to connect their computer to the Internet. This shows that the use of Internet networks is very important, whether they used it for work and assignment purposes, or for the access to social media accounts, such as Instagram, Facebook and Twitter. However, in this wide use of this computer network, the privacy of computer users is in danger. This is because some of the computer users do not install security system in their computer. This problem will allow the hackers to hack and commit the network attacks. This is very dangerous, especially to the important organizations because hackers can disable the online system in the company, steal confidential information and subsequently steal company money through online without being aware of any one. The attacks that can be made includes denial-of-Service attack, DNS spoofing attack and phishing attack. The goal of this study is to apply anti-phishing tools in preventing the network security attack in an organization. In this study, phishing attacks have been studied thoroughly. After a study has been made, machine learning method is used to prevent the phishing attack. Besides, the study also shows that phishing attack is always related to the spam attack, where there might be attached phishing link in the spam message. This spam message includes the email and the SMS message that received by the user. There are several algorithms that can be used in the machine learning method to prevent the both attacks. The Naïve Bayes algorithm, Decision Tree algorithm and Support Vector Machine algorithm has been used to prevent the spam attack, as well as the phishing attack. The study of this algorithm is made thoroughly and the methods in implementing this algorithm have been discussed in detail. The experiment is conducted for the datasets that obtained by using machine learning method. The results are obtained, showing the performance of machine learning method on each dataset.

TABLE OF CONTENT

DECLARATION

TITLE PAGE

ACKNOWLEDGEMENTS **ii**

ABSTRAK **iii**

ABSTRACT **iv**

TABLE OF CONTENT **v-vi**

LIST OF TABLES **vii**

LIST OF FIGURES **viii**

LIST OF ABBREVIATIONS **ix**

CHAPTER 1 INTRODUCTION **1**

1.1 Background of study 1

1.2 Problem statement 2

1.3 Research goal and objective 3

1.4 Scope of research 3

1.5 Significance of research 4

1.6 Report organisation 4

CHAPTER 2 LITERATURE REVIEW **5**

2.1 Introduction 5

2.2 Phishing attack 5

2.3 Existing Machine Learning method 6

2.3.1 Naïve Bayes algorithm 6

2.3.2	Decision tree algorithm	8
2.3.3	Support Vector Machine algorithm	9
2.3.4	Comparison among three existing algorithm	11
2.4	Technologies	14
2.5	Conclusion	15
CHAPTER 3 METHODOLOGY		16
3.1	Introduction	16
3.2	Methodology	16
3.2.1	Literature study	16
3.2.2	Data collection	17
3.2.3	Classification	18
3.2.4	Performance measurement	21
3.3	Hardware and software	22
3.4	Gantt chart	23
3.5	Implementation	26
3.6	Conclusion	28
CHAPTER 4 TESTING AND RESULT DISCUSSION		29
4.1	Introduction	29
4.2	Results	29
4.2.1	Dataset 1: The email messages	29
4.2.2	Dataset 2: The SMS messages	31
4.3	Discussion	34
4.4	Conclusion	38

CHAPTER 5 CONCLUSION	39
5.1 Concluding remarks	39
5.2 Research constraints and challenges	40
5.3 Future work	40
REFERENCES	41

LIST OF TABLES

Table 1.1	Tabulation of the problem	3
Table 2.1	The performance measurement	7
Table 2.2	Tabulation of the SVM algorithm's attribute and the significance of the attributes	10
Table 2.3	Tabulation of the three algorithms to prevent the phishing attack.	12
Table 2.4	Tabulation of the types of programming languages.	14
Table 3.1	The hardware requirements and specification.	22
Table 3.2	The software requirements and specification.	23
Table 4.1	Tabulation of the result of dataset 1	30
Table 4.2	Tabulation of the result of dataset 2	32
Table 4.3	Tabulation of the result based on the previous research (Metsis, Androutsopoulos, & Paliouras, 2006)	36
Table 4.4	Tabulation of the result based on the previous research (Almedia, Hidalgo, & Yamakami, 2011)	36

LIST OF FIGURES

Figure 3.1	The steps in the research methodology.	17
Figure 3.2	The Gantt chart from phase 1 to phase 2.	24
Figure 3.3	The Gantt chart from phase 3 to phase 4.	24
Figure 3.4	The Gantt chart from phase 4 to phase 5.	25
Figure 3.5	The summary of the implementation process.	27
Figure 4.1	Graph of percentage of correctly classified instances based on Dataset 1.	31
Figure 4.2	Graph of percentage of correctly classified instances based on Dataset 2.	33
Figure 4.3	Graph of time taken for the classification of dataset using Naïve Bayes algorithm.	34
Figure 4.4	Graph of time taken for the classification of dataset using Decision Tree algorithm.	35
Figure 4.5	Graph of time taken for the classification of dataset using Support Vector Machine algorithm.	35

LIST OF ABBREVIATIONS

DNS	Domain Name System
URL	Uniform Resource Locator
SVM	Support Vector Machine
ARFF	Attribute Relation File Format

CHAPTER 1

INTRODUCTION

1.1 Background of study

Network security is the most important and critical issues that need to be considered and emphasized in the network, especially in an organization, such as offices, banks and clinics. Basically, network security is the authorization, commonly by using a username and password, which inhibit and monitor the unauthorized access and all the administrator event in the network (Pawar & Anuradha, 2015).

It is important for the organization to maintain their security network to ensure the privacy and confidentiality of their employer data, as well as their employee data. This will ensure the data, especially the sensitive data, such as the employee information details, can be stored in the server safely. For example, for us to access the online banking, we need to have an authentication to access our account. This can be done by providing the username and password in the login page of the online banking. Authentication is needed in this scenario so that our sensitive data would not be exposed to the unauthorized user or the hacker.

Although there is implementation of the network security in an organization, but still there is network attack happened. The network attack that usually happen includes phishing, denial-of-Service attack and Domain Name System (DNS) spoofing. This attack will contribute to the financially and privacy loss to the victims. For example, when the hacker attacks sensitive information while the user using their online banking account, the attacker will use this information to retrieve back the victim's account and then steal their money inside the account. This also can be applied to the office organization, whereby the hacker will gain the sensitive data and use it to commit online crimes, such as stealing the office's money and the data of their employer over the Internet.

Phishing is one of the network security attack, which is the derivational of word 'Fishing' by replacing the 'F' with 'Ph'. Phishing is the act of imitate the genuine websites to collect the sensitive information from the victim and use it for committing crimes, such as illegal financial gain (Kaur & Kaur, 2015). This attack typically starts when the hacker sends an email that seems original to the victim and persuade them to update and verify their information by clicking the Uniform Resource Locator (URL) link in the email (Mohammad, Thabtah, & McCluskey, 2015). Usually, the phishing email will redirect the user to the infected website and asking them to provide their particular information, such as their personal details and bank account information, which will be used to hack the information whatever the user enter (Suganya, 2016). The phishing attack is always related to the spamming email that received by the victim. Those spam emails are also vulnerable to the phishing attack because some of the spam email may contain the link that will redirect the victim to the phishing websites.

The phishing attack can be prevented using the machine learning method. According to Marsland (2015), machine learning is the modification or adaptation of the computer actions so that we can get the more accurate actions in the end. Besides, machine learning is also considered as computational complex since it will lead to produce algorithm. Based on the machine learning method, the prevention of phishing attack can be classified to several algorithm, which includes Decision Tree algorithm, Naïve Bayes algorithm and Support Vector Machine (SVM) algorithm (Smadi, Aslam, Zhang, Alasem, & Hossain, 2016). All the algorithms stated are used in classifying the spam email and the SMS messages datasets. This will show the performance of each of the algorithm, in terms of their accuracy.

1.2 Problem statement

After doing some research, there are some methodologies of overcoming the attack that has been found, which is already exists. Each of them has their own advantages and disadvantages respectively. The summary of the problems is tabulated in the Table 1.1;

Table 1.1: Tabulation of the problem

No.	Problem	Description	Effect
1	Classification of the phishing result is not accurate (Smadi et al., 2016).	The result that obtained after the test, which is the true positive and the false positive does not be considered.	The result that obtained might not be correct and accurate.
2	The level of performance of the method (Smadi et al., 2016).	The performance level of the method does not be considered when the method was used.	The result of the test might take longer time to be obtain.

1.3 Research goal and objective

The goal of this research is to detect phishing attack. The objective of this study is stated as below;

- i. To study the issues of phishing attack.
- ii. To use machine learning method, which is Naïve Bayes algorithm, Decision Tree algorithm and Support Vector Machine algorithm in detection of phishing attack.
- iii. To evaluate the performance of machine learning method in detection of phishing attack.

1.4 Scope of the research

The scope of this research is listed as follow;

- i. The research is focus on the method on how to overcome the network security attack, which is machine learning method.
- ii. The network security attack will be observed thoroughly.
- iii. The dataset that will be use in this algorithm is the email message dataset and SMS message dataset, which contains spam and legitimate (ham) messages.

REFERENCES

- A. Yasin and A. Abuhasan, “An Intelligent Classification Model For Phishing Email Detection,” *International Journal of Network Security & Its Applications (IJNSA)*, vol. 8, no. 4, 2016.
- Aized Amin Soofi and Arshad Awan, “Classification Techniques in Machine Learning: Applications and Issues,” *Journal of Basic & Applied Sciences*, vol. 13, pp. 459–465, 2017.
- Akansha, P., & Meenakshi, E. (2017). Detection of Phishing Websites Using Data Mining Techniques, 2(12), 1468–1472.
- Ali, W. (2017). Spam Detection using WEKA. Retrieved from <https://github.com/waleedalinizami/Spam-Detection-Using-Weka>
- Almeida, T. A., Hidalgo, J. M. G., & Yamakami, A. (2011). Contributions to the study of SMS spam filtering. *Proceedings of the 11th ACM Symposium on Document Engineering - DocEng '11*, 259. <https://doi.org/10.1145/2034691.2034742>
- Almeida, T. A., & José María Gómez Hidalgo. (2011). SMS Spam Collection v.1. Retrieved from <http://dcomp.sor.ufscar.br/talmeida/smspamcollection/>
- Bird, M. (2015). The Pros and Cons of using Java. Retrieved March 22, 2018, from <http://www.digitalrise.biz/software/the-pros-and-cons-of-using-java/>
- Gupta, P. (2017). Cross-Validation in Machine Learning. Retrieved from <https://towardsdatascience.com/cross-validation-in-machine-learning-72924a69872f>
- H. Byun and S.-W. Lee, “Applications of Support Vector Machines for Pattern Recognition: A Survey,” in *Pattern Recognition with Support Vector Machines*, 2002, pp. 213–236.
- Jasmina Novaković, Perica Strbac, and Dusan Bulatović, “Toward optimal feature selection using ranking methods and classification algorithms,” *Yugoslav Journal of Operations Research*, vol. 21, no. 1, pp. 191–135, 2011.
- Jayesh Bapu Ahire. (2016). Introduction to Word Vector, 1–11. Retrieved from <https://medium.com/@jayeshbahire/introduction-to-word-vectors-ea1d4e4b84bf>

- K. Kim, "A hybrid classification algorithm by subspace partitioning through semi-supervised decision tree," *Pattern Recognition*, vol. 60, pp. 157–163, Dec. 2016.
- Kaur, S., & Kaur, A. (2015). Detection of Phishing Webpages using Weights computed through Genetic Algorithm. In *2015 IEEE 3rd International Conference on MOOCs, Innovation and Technology in Education (MITE)* (pp. 331–336). Amritsar: IEEE.
- Kozak, J., & Boryczka, U. (2016). Collective data mining in the ant colony decision tree approach. *Information Sciences*, 372, 126–147. <https://doi.org/10.1016/j.ins.2016.08.051>
- Kumar, N., & Chaudhary, P. (2017). Mobile Phishing Detection using Naive Bayesian Algorithm, *17*(7), 142–147.
- Lysis. (2017). Pros and Cons of Using C# as Your Backend Programming Language. Retrieved March 22, 2018, from <https://agilites.com/pros-and-cons-of-using-c-as-your-backend-programming-language.html>
- Mack, D. (2018). How to pick the best learning rate for your machine learning project. Retrieved from <https://medium.freecodecamp.org/how-to-pick-the-best-learning-rate-for-your-machine-learning-project-9c28865039a8>
- Marsland, S. (2015). *Machine Learning: An Algorithmic Perspective* (2nd ed.). Boca Raton: Taylor & Francis Group.
- Metsis, V., Androutsopoulos, I., & Paliouras, G. (2006). Spam filtering using Naive Bayes-Which Naive Bayes? <https://doi.org/10.3109/02841866309134119>
- Mohammad, R. M., Thabtah, F., & McCluskey, L. (2015). Tutorial and critical analysis of phishing websites methods. *Computer Science Review*, 17, 1–24. <https://doi.org/10.1016/j.cosrev.2015.04.001>
- P. Ghamisi, J. Plaza, Y. Chen, J. Li, and A. J. Plaza, "Advanced Spectral Classifiers for Hyperspectral Images: A review," *IEEE Geoscience and Remote Sensing Magazine*, vol. 5, no. 1, pp. 8–32, 2017.
- Patil, P., Rane, R., & Bhalekar, M. (2017). Detecting spam and phishing mails using SVM and obfuscation URL detection algorithm. *Proceedings of the International Conference on Inventive Systems and Control, ICISC 2017*, 1–4. <https://doi.org/10.1109/ICISC.2017.8068633>

- Patil, T. R. (2013). Performance Analysis of Naive Bayes and J48 Classification Algorithm for Data Classification. *International Journal Of Computer Science And Applications*, ISSN: 0974-1011, 6(2), 256–261. <https://doi.org/ISSN: 0974-1011>
- Pawar, M. V., & Anuradha, J. (2015). Network security and types of attacks in network. *Procedia Computer Science*, 48(C), 503–506. <https://doi.org/10.1016/j.procs.2015.04.126>
- Rathod, S. B., & Pattewar, T. M. (2015). Content Based Spam Detection in Email using Bayesian Classifier. In *Conference: 2015 International Conference on Communications and Signal Processing (ICCSP)* (pp. 1257–1261). <https://doi.org/10.1109/ICCSP.2015.7322709>
- Ray, S. (2017). Understanding Support Vector Machine algorithm from examples (along with code). Retrieved March 22, 2018, from <https://www.analyticsvidhya.com/blog/2017/09/understaing-support-vector-machine-example-code/>
- S.Archana and K.Elangovan, “Survey of Classification Techniques in Data Mining,” *International Journal of Computer Science and Mobile Applications*, vol. 2, no. 2, pp. 65–71, 2014.
- Sahouane, A. (2016). The pros and cons of Python. Retrieved March 22, 2018, from <https://www.supinfo.com/articles/single/3425-the-pros-and-cons-of-python>
- Sao, P., & Prashanthi, P. K. (2015). E-mail Spam Classification Using Naïve Bayesian Classifier. *International Journal of Advanced Research in Computer Engineering & Technology (IJARCET)*, 4(6), 2792–2796.
- Saxena, R. (2017a). How Decision Tree Algorithm Works. Retrieved March 22, 2018, from <http://dataaspirant.com/2017/01/30/how-decision-tree-algorithm-works/>
- Saxena, R. (2017b). SVM Classifier, Introduction to Support Vector Machine Algorithm. Retrieved April 5, 2018, from <http://dataaspirant.com/2017/01/13/support-vector-machine-algorithm/>
- Smadi, S., Aslam, N., Zhang, L., Alasem, R., & Hossain, M. A. (2016). Detection of phishing emails using data mining algorithms. In *SKIMA 2015 - 9th International Conference on Software, Knowledge, Information Management and Applications* (pp. 1–8). <https://doi.org/10.1109/SKIMA.2015.7399985>

- Suganya, V. (2016). A Review on Phishing Attacks and Various Anti Phishing Techniques. *International Journal of Computer Applications*, 139(1), 975–8887.
- Waldron, M. (2015). Naive Bayes for Dummies; A Simple Explanation. Retrieved March 22, 2018, from <http://blog.aylien.com/naive-bayes-for-dummies-a-simple-explanation/>
- Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining: Practical Machine Learning Tools and Techniques* (Fourth). Cambridge: Elsevier Inc.
- Yang, X., Yan, L., Yang, B., & Li, Y. (2017). Phishing Website Detection Using C4 . 5 Decision Tree, (Itme), 119–124.
- Yitagesu, M. E., & Tijare, P. M. (2016). Email Classification using Classification Method, 32(3), 142–145.